# SOCR: A Handwritten Data Form Producing and Reading System

**Hanchuan Peng and Qiang Gan**

Department of Biomedical Engineering,
Southeast University,
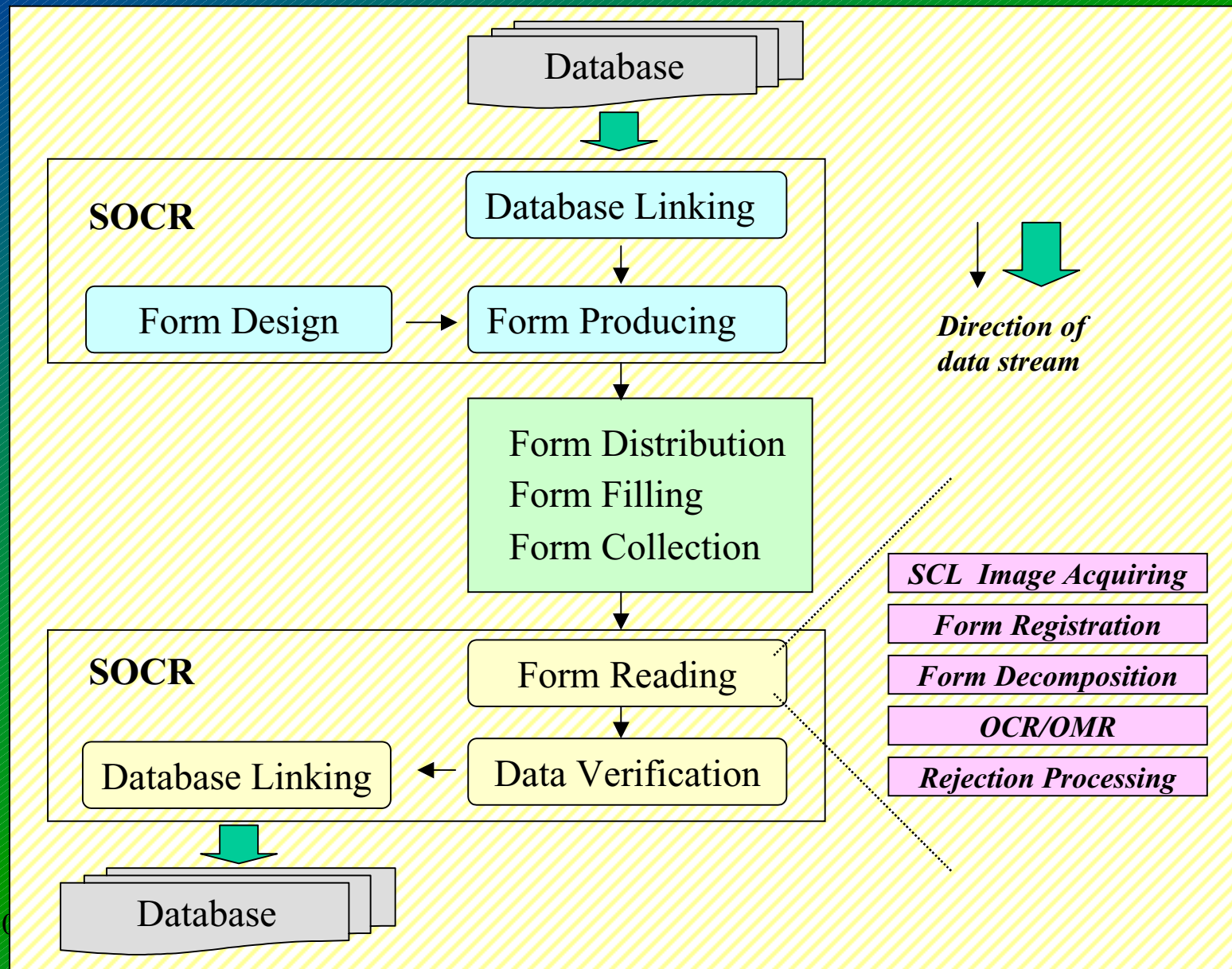Nanjing, 210096, China.
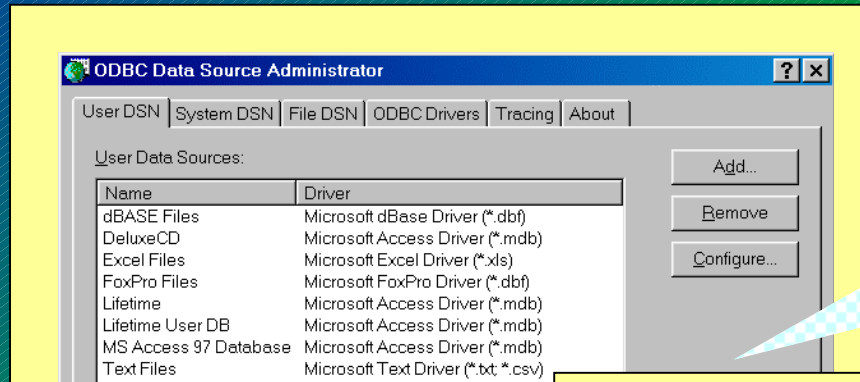Email: phc@seu.edu.cn

2000-Jan-10

# Abstract

*Document analysis and recognition software, especially form reader, is greatly required in office automation. This paper serves as a primer to the SOCR (version 1.03), which is a software package for handwritten data form producing and reading. This package integrates techniques of database linking, form producing, form reading, and data verification. With SOCR, data forms can be easily created and the form data, including loosely constrained handwritten numerals and symbols, can be read into databases at high accuracy and high speed, and in a human-compatible manner. This software package has been applied to producing and reading handwritten student score forms in several universities and tax forms in several cities.*

2000-Jan-10

# Main Structure of the System

Database

↓

**SOCR**

Database Linking

Form Design → Form Producing

↓

Form Distribution
Form Filling
Form Collection

↓

**SOCR**

Form Reading

↓

Database Linking ← Data Verification

↓

Database

*Direction of data stream*

*SCL Image Acquiring*
*Form Registration*
*Form Decomposition*
*OCR/OMR*
*Rejection Processing*

# Database Linking and Form Producing



**SOCR can access existing databases and enable seamless usage of data forms in the databases and form data in paper documents.**
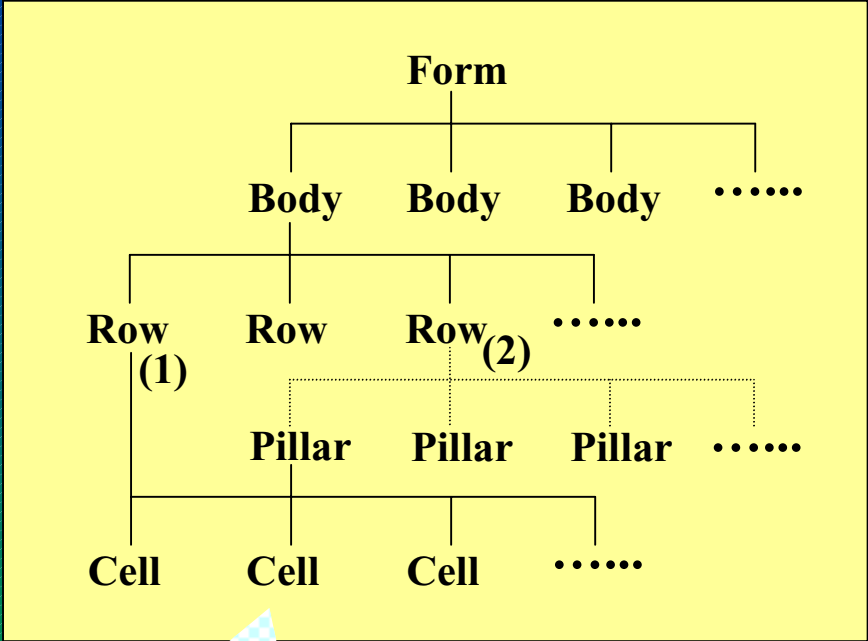
**ODBC and DAO are used to link various databases**

**Functions are provided to merge data fields (in database) and to link to specified fields of form templates.**

**Form Producing Wizard**

# Form Data Structure

**Form**

**Body** **Body** **Body** ‥‥‥

**Row (1)** **Row** **Row (2)** ‥‥‥

**Pillar** **Pillar** **Pillar** ‥‥‥

**Cell** **Cell** **Cell** ‥‥‥

*The hierarchy tree, where the level of "Pillar" is optional, depending on applications.*

*An example of form without the "Pillar" level.*

B1R1C1

B2R2C1

B2R3C2

B3R1C1    B3R1C3

B3R2C2

B1R2P1C1    B1R2P2C1

B1R2P1R2    B1R2P2C2

B2R3P1C1  B2R3P2C1LT

B2R3P1C2    B2R3P2C1RT

B3R1P1C1RT

B3R1P1C1LB

B3R2C2 or
B3R2P2C1

*An example of form without the "Pillar" level.*

2000-Jan-10

# Sample Form

## 东南大学学生成绩登记表

| 课程名称 | 美术及艺术 | 课程代号 | 000001 | 教师姓名 | 教一 | 教师代号 | 000001 | 考试类别 | | 学分 | 3·5 |
|---|---|---|---|---|---|---|---|---|---|---|---|

（注： 考试类别编号为，1-考试，2-考查，3-补考，4-重修）

| 序号 | 学号 | 姓名 | 分数 | 更正分 | 序号 | 学号 | 姓名 | 分数 | 更正分 | 序号 | 学号 | 姓名 | 分数 | 更正分 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 01 | 11191101 | 学一 | 60 | 75 | 16 | 11191116 | 学十六 | 100 | | 31 | 11191131 | 学三一 | ■ | 53 |
| 02 | 11191102 | 学二 | 61 | 76 | 17 | 11191117 | 学十七 | 99 | | 32 | 11191132 | 学三二 | ■ | 60 |
| 03 | 11191103 | 学三 | 62 | 77 | 18 | 11191118 | 学十八 | 98 | | 33 | 11191133 | 学三三 | 71 | |
| 04 | 11191104 | 学四 | 63 | 78 | 19 | 11191119 | 学十九 | 97 | | 34 | 11191134 | 学三四 | 72 | |
| 05 | 11191105 | | 64 | 79 | 20 | 11191120 | 学二十 | 96 | | 35 | 11191135 | 学三五 | 73 | |
| | | | | 80 | 21 | 11191121 | 学二一 | 95 | | 36 | 11191136 | 学三六 | 74 | |
| | | | | | 22 | 11191122 | 学二二 | ■ | 94 | 37 | 11191137 | 学三七 | 75 | |
| | | | | | 23 | 11191123 | 学二三 | 93 | | 38 | 11191138 | 学三八 | 76 | |
| | | | | | 24 | 11191124 | 学二四 | 92 | | 39 | 11191139 | 学三九 | 77 | |
| | | | 69 | 84 | 25 | 11191125 | 学二五 | 91 | | 40 | 11191 | | | |
| 11 | 11191111 | 学十一 | 70 | 85 | 26 | 11191126 | 学二六 | 90 | | 4 | | | | |
| 12 | 11191112 | 学十二 | 71 | 86 | 27 | 11191127 | 学二七 | 89 | | | | | | |
| 13 | 11191113 | 学十三 | 72 | 87 | 28 | 11191128 | 学二八 | 88 | | | | | | |
| 14 | 11191114 | 学十四 | 73 | 88 | 29 | 11191129 | 学二九 | 87 | | 44 | | | | |
| 15 | 11191115 | 学十五 | 74 | 87 | 30 | 11191130 | 学三十 | 86 | | 45 | 11191145 | 学四五 | | 05 |

书写要求：
（1）黑钢笔书写，笔划横平竖直，书写位置居中，勿与边框相连，勿出现断笔。
（2）数字1、2、3、4、5、7不带圈，0、6、8、9圈要闭合。
（3）数字4上部要有较大开口。
（4）误写字符请涂黑。

规范书写： 0123456789■

### 更正计录

| 序号 | 签名 | 盖章 |
|---|---|---|
| | | |
| | | |
| | | |

教师签名　　　　　教学系主任签名

年　月　日　　　年　月　日

*The form format can be changed by the form producing module*

*Form registration is automatically performed. A wizard is provided.*

# Handwritten Character Recognition

- The OCR/OMR module of SOCR can attain about **99.5%** recognition accuracy for loosely constrained (handprinted) handwritten numerals and symbols.

- "Loosely constrained" means that all the characters should not be indeterminate for classification. The constraints imposed on handwritten numerals include:
  - '0', '6', '8', '9' have closed circles,
  - '1', '2', '3', '4', '5', '7' do not have circles,
  - '4' does not close the upper part.

- The recognition engine of SOCR is based on a hybrid model, which combines both traditional methods and neural recognizers. Doubtful characters are sent to the rejection processing module.

2000-Jan-10

# Rejection Processing and Data Verification



*Manual tool is provided for rejection correction. All the corrected characters can be distributed back to their original forms automatically.*

*Different data verification strategies are used to examine the (logical) correctness of recognized results. A pair examination wizard is provided.*

2000-Jan-10